# NLIDB Converter for Customer Relationship Index

G. O. Wijayasekara

149236 L

Dissertation submitted to the Faculty of Information Technology, University of

Moratuwa, Sri Lanka for the partial fulfillment of the requirements of the Degree of

Master of Science in Information Technology.

May 2017

**Declaration**

I declare that this thesis is my own work and has not been submitted in any form for another degree or diploma at any university or other institution of tertiary education. Information derived from the published or unpublished work of others has been acknowledged in the text and a list of references is given.


**Gayanee Wijayasekara**

Name of Student

…………………………………

Signature of Student


…………………………………

Date


**S.C. Premaratne**

Name of Supervisor

…………………………………

Signature of Supervisor


.................................................

Date

**Acknowledgement**

I take this opportunity to express my heartfelt gratitude to my supervisor, Mr. S. C. Premaratne, Senior Lecturer at the University of Moratuwa, Sri Lanka, for guiding me throughout the research project without whose support and direction this would not have been possible.

And a special heartfelt thanks goes out to my loving family for being the pillars of strength that they are at all times.

I would also like to thank my colleagues for the valuable feedback and support they extended during the course of the research.

**Abstract**

Information plays a significant role in our day to day life. Information affects how we look at things and make decisions. Therefore, accurate information must be readily available. With technology reaching new heights, in modern day systems and applications, the main source for information storage is an underlying database system. Retrieving information from a database management system requires a specific expert skill set which predominantly includes knowledge on Structured Query Language and domain specific knowledge. In recent times there is a rising demand for non-expert users to be able to directly extract information from an underlying database management system. For industry specific Customer Relationship Management applications like the Customer Relationship Index used in large organizations like the Brandix Group this is seen as both a challenge and an opportunity to be explored. Implementation of a Natural Language Interface to the underlying database of the Customer Relationship Index allows a typical user to retrieve the required information from the underlying database using natural language like English without prior programming or technical knowledge. For above reasons the NLIDB converter for the Customer Relationship Index is introduced. This application takes the users query to be run against the underlying database in natural language and returns the corresponding T-SQL statement which can then be run against the database to extract results. The NLIDB converter has shown that it can successfully and efficiently convert questions in natural language to the corresponding T-SQL statement with an accuracy rate of more than 80% for two types of output T-SQL statement formats for the Customer Relationship Index database.

**Contents**

**List of Figures**

x

**List of Tables**